

## Association for Information Systems AIS Electronic Library (AISeL)

---

PACIS 2006 Proceedings

Pacific Asia Conference on Information Systems  
(PACIS)

---

2006

# A Close Look at Privacy Preserving Data Mining Methods

Xiaodan Wu

*Hebei Univ. of Technology, xwu@hebut.edu.cn*

Yunfeng Wang

*Hebei Univ. of Technology, ywang@hebut.edu.cn*

Chao-Hsien Chu

*Singapore Management University, chuch@smu.edu.sg*

Fengli Liu

*Hebei Univ. of Technology, Liufengli312@163.com*

Ping Chen

*Hebei Univ. of Technology, Cplb203@yahoo.com.cn*

*See next page for additional authors*

Follow this and additional works at: <http://aisel.aisnet.org/pacis2006>

---

### Recommended Citation

Wu, Xiaodan; Wang, Yunfeng; Chu, Chao-Hsien; Liu, Fengli; Chen, Ping; and Yue, Dianmin, "A Close Look at Privacy Preserving Data Mining Methods" (2006). *PACIS 2006 Proceedings*. 32.

<http://aisel.aisnet.org/pacis2006/32>

This material is brought to you by the Pacific Asia Conference on Information Systems (PACIS) at AIS Electronic Library (AISeL). It has been accepted for inclusion in PACIS 2006 Proceedings by an authorized administrator of AIS Electronic Library (AISeL). For more information, please contact [elibrary@aisnet.org](mailto:elibrary@aisnet.org).

---

**Authors**

Xiaodan Wu, Yunfeng Wang, Chao-Hsien Chu, Fengli Liu, Ping Chen, and Dianmin Yue

## A Close Look at Privacy Preserving Data Mining Methods

Xiaodan Wu  
Hebei Univ. of Technology  
Tianjin 300130, P.R.C.  
xwu@hebut.edu.cn

Yunfeng Wang  
Hebei Univ. of Technology  
Tianjin 300130, P.R.C.  
ywang@hebut.edu.cn

Chao-Hsien Chu  
Singapore Management  
Univ., Singapore 178902  
chuch@smu.edu.sg

Fengli Liu  
Hebei Univ. of Technology  
Tianjin 300130, P.R.C.  
Liufengli312@163.com

Ping Chen  
Hebei Univ. of Technology  
Tianjin 300130, P.R.C.  
Cplb203@yahoo.com.cn

Dianmin Yue  
Hebei Univ. of Technology  
Tianjin 300130, P.R.C.

### Abstract

*Recent advances in information, communications, data mining, and security technologies have gave rise to a new era of research, known as privacy preserving data mining (PPDM). Several data mining algorithms, incorporating privacy preserving mechanisms, have been developed that allow one to extract relevant knowledge from large amount of data, while hide sensitive data or information from disclosure or inference. PPDM is a new attempt; thus, several research questions have often being asked. For instance: (1) how to measure the performance of these algorithms? (2) how effective of these algorithms in terms of privacy preserving? (3) will they impact the accuracy of data mining results? And (4) which one can better protect sensitive information? To help answer these questions, we conduct an extensive review on literature. We present a classification scheme, adopted from early studies, to guide the review process. Finally, we share directions for future research.*

**Keywords:** Privacy preservation, data mining, knowledge discovery, data perturbation, secure multiparty computation

### 1. Introduction

Increasing network complexity, affording greater access, sharing information and a growing emphasis on the Internet have made information security and privacy a major concern for individuals and organizations. Data mining is a well-known technology for automatically and intelligently extracting knowledge from large amount of data. Such a process, however, can also disclosure sensitive information about individuals compromising the individual's right to privacy. Moreover, data mining techniques can reveal critical information about business transactions, compromising the free competition in a business setting (Bertino et. al. 2005). Privacy preserving data mining (PPDM) is a new era of research in data mining. Its ultimate goal is to develop efficient algorithms that allow one to extract relevant knowledge from large amount of data, while prevent sensitive information from disclosure or inference.

PPDM research usually takes one of the three philosophical approaches: (1) data hiding, in which sensitive raw data like identifiers, name, addresses, etc. were altered, blocked, or trimmed out from the original database, in order for the users of the data not to be able to compromise another person's privacy; (2) rule hiding, in which sensitive knowledge extracted from the data mining process be excluded for use, because confidential information

may be derived from the released knowledge. This problem is also commonly called the “database inference problem;” and (3) secure multiparty computation (SMC), where distributed data are encrypted before released or shared for computations; thus, no party knows anything except its own inputs and the results.

PPDM is a fast growing research area. Given the number of different algorithms have been developed over the last years, there is an emerging need of synthesizing literature to understand the nature of problem, identify potential research issues, standardize new research area, and evaluate the relative performance of different approaches (Verykios et al. 2004; Bertino et al. 2005). The main purpose of this study is to review the state-of-the-art in current PPDM research in order to better understand existing algorithms, answer research questions and move forward the field of research.

## 2. Classification Framework for PPDM

In this paper, we propose to consolidate and simplify the taxonomy brought by Bertino et al. (2005). We propose to reduce the PPDM taxonomy into four levels: data distribution, purposes of hiding, data mining algorithms, and privacy preserving techniques (see Figure 1).

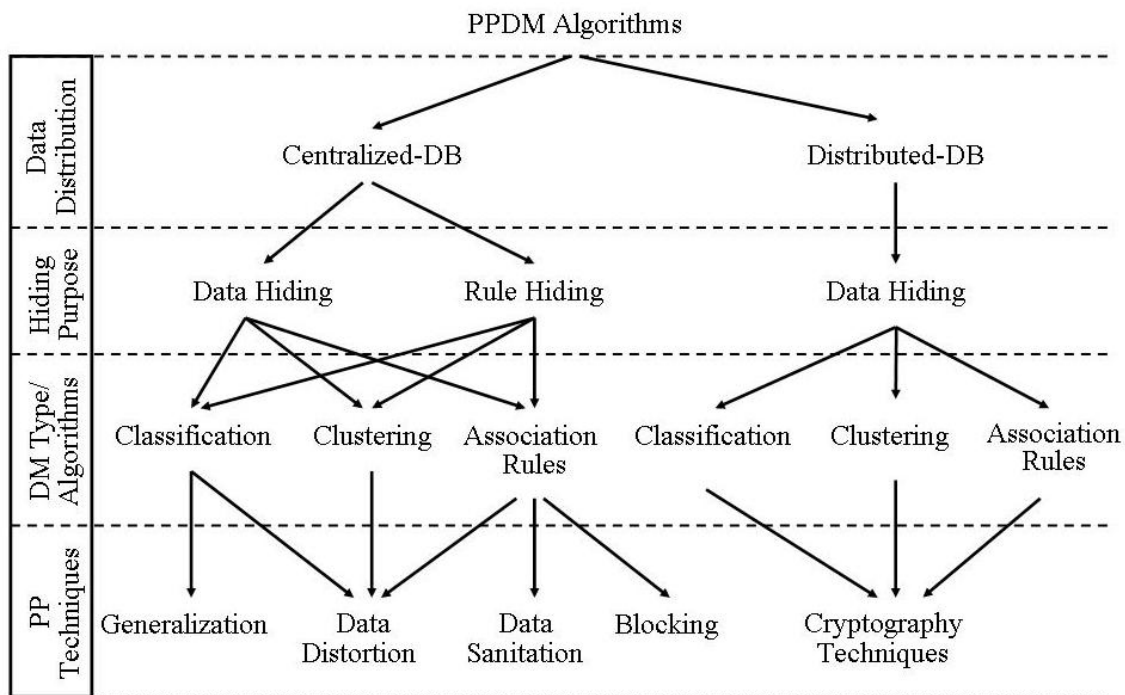


Figure 1: The Taxonomy of PPDM algorithms

### 2.1 Data Distribution

The PPDM algorithms can be first divided into two major categories, centralized and distributed data, based on the distribution of data. In a centralized database (DB) environment, data are all stored in a single database; while, in a distributed database environment, data are stored in different databases. Earlier research has been predominately focused on dealing with privacy preservation in a centralized DB (Du and Zhan 2003; Evfimievski et al. 2003, 2004; Islam and Brankovic 2004; Natwichail, et al. 2005; Oliveira, et al. 2002, 2003a, 2003b,

2003c, 2004a, 2004b; Rizvi and Haritsa 2002; Saygin et al. 2001, 2002; Verkios, et al. 2003; Wang et al. 2004; Xia, et al. 2004). The difficulties of applying PPDM algorithms to a distributed DB can be attributed to two reasons: first, the data owners have privacy concerns so they may not willing to release their own data for others; second, even if they are willing to share data for data mining, the communication cost between the sites is too expensive. In today's global digital environment, most data are often stored in different sites, thus, more attention and research should be focused on distributed PPDM algorithms.

## **2.2 Hiding Purposes**

The PPDM algorithms can be further classified into two types, data hiding and rule hiding, according to the purposes of hiding. Data hiding refers to the cases where the sensitive data from original database like identity, name, and address that can be linked, directly or indirectly, to an individual person are hided. In contrast, in rule hiding, we remove the sensitive knowledge derived from original database after applying data mining algorithms. Majority of the PPDM algorithms used data hiding techniques. This is especially true in a distributed database environment (Du and Zhan 2002; Kantarcioglu and Clifton, 2002, 2003; Klusch et al. 2003; Lindell and Pinkas 2000; Merugu and Ghosh 2003; Vaidya and Clifton 2002, 2003, 2005; Verkios, et al. 2003; Yang et al. 2006; Zhan et al. 2005), as the techniques can be used to prevent individual information from being discovered by other parties in the joint computational process. Most PPDM algorithms hide sensitive patterns by modifying data (Du and Zhan 2003; Oliveira, et al. 2003c, 2004a; Xia, et al. 2004; Du and Zhan 2002; Evfimievski et al. 2003, 2004; Kantarcioglu and Clifton 2002, 2003; Islam and Brankovic 2004; Klusch et al. 2003; Lindell and Pinkas 2000; Merugu and Ghosh 2003; Rizvi and Haritsa 2002; Vaidya and Clifton 2002, 2003, 2005; Verkios et al. 2003; Wang et al. 2004; Yang et al. 2006; Zhan et al. 2005). Also, at present, the rule hiding techniques is only being adopted by association rule mining for centralized DB (Oliveira et al. 2002, 2003a, 2003b, 2004b; Verkios et al. 2003; Saygin et al. 2001, 2002).

## **2.3 Data Mining Tasks/Algorithms**

Currently, the PPDM algorithms are mainly used on the tasks of classification, association rule and clustering. Association analysis involves the discovery of associated rules, showing attribute value and conditions that occur frequently in a given set of data. Classification is the process of finding a set of models (or functions) that describe and distinguish data classes or concepts, for the purpose of being able to use the model to predict the class of objects whose class label is unknown. Clustering Analysis concerns the problem of decomposing or partitioning a data set (usually multivariate) into groups so that the points in one group are similar to each other and are as different as possible from the points in other groups. A majority of the PPDM algorithms used association rule method for mining data (Evfimievski et al. 2003, 2004; Oliveira et al. 2002, 2003a, 2003b, 2004b; Rizvi and Haritsa 2002; Saygin et al. 2001, 2002; Verkios et al., 2003; Xia et al. 2004; Kantarcioglu and Clifton 2002, 2003; Vaidya and Clifton 2005; Veloso et al. 2003), followed by classification (Du and Zhan 2003; Islam and Brankovic 2004; Natwichail, et al. 2005; Wang et al. 2004; Du and Zhan 2002; Kantarcioglu and Clifton 2003; Lindell and Pinkas 2000; Vaidya and Clifton 2005; Yang et al. 2006;), and then clustering (Oliveira et al. 2003c, 2004a; Vaidya and Clifton 2003; Klusch et al. 2003; Merugu and Ghosh 2003).

## **2.4 Privacy Preservation Technique**

Four techniques – sanitation, blocking, distort, and generalization -- have been used to hide data items for a centralized data distribution. Data sanitation is to remove or modify items in a database to reduce the support of some frequently used itemsets such that sensitive patterns cannot be mined. The blocking approach replaces certain attributes of the data with a question mark. In this regard, the minimum support and confidence level will be altered into a minimum interval. As long as the support and/or the confidence of a sensitive rule lie below the middle in these two ranges, the confidentiality of data is expected to be protected. Data distort protects privacy for individual data records through modification of its original data, in which the original distribution of the data is reconstructed from the randomized data. These techniques aim to design distortion methods after which the true value of any individual record is difficult to ascertain, but “global” properties of the data remain largely unchanged. Generalization transforms and replaces each record value with a corresponding generalized value.

The privacy preservation technique used in a distributed database is mainly based on cryptography techniques. SMC algorithms deal with computing any function on any input, in a distributed network where each participant holds one of the inputs, while ensuring that no more information is revealed to a participant in the computation than can be inferred from that participant’s input and output. Data distort is the most popular method used in hiding data (Du and Zhan 2003; Evfimievski et al. 2003, 2004; Islam and Brankovic 2004; Oliveira et al., 2003c, 2004a; Rizvi and Haritsa 2002; Xia et al. 2004), followed by data sanitation (Oliveira et al. 2002, 2003a, 2003b, 2004b; Saygin et al. 2002; Verkios et al. 2003) and generalization (Natwichail et al. 2005; Wang et al. 2004). If one wants to obtain data mining results from different data sources, then the only method can be used is a cryptography technique (Du and Zhan 2002; Kantarcioglu and Clifton 2002, 2003; Klusch et al. 2003; Lindell and Pinkas 2000; Merugu and Ghosh 2003; Vaidya and Clifton 2002, 2003, 2005; Veloso et al. 2003; Yang et al. 2006; Zhan et al. 2005). Since the parties who use SMC operators cannot reveal anything from others except final results, it can have benefits of both accuracy of data mining results and the privacy of the database.

### ***3. Suggestions for Future Work***

First, current studies tend to use different terminology to describe similar or related practice. For instance, people have used data modification, data perturbation, data sanitation, data hiding, and preprocessing as possible methods for preserving privacy; however, all are in fact related to the use of some types of technique to modify original data so that private data and knowledge remain private even after the mining process. Lacking a common language for discussions will cause misunderstanding and slow down the research breakthrough. Therefore, there is an emerging need of standardizing the terminology and PPDM practice.

Second, most prior PPDM algorithms were developed for use with data stored in a centralized database. However, in today’s global digital environment, data is often stored in different sites. With recent advances in information and communication technologies, the distributed PPDM methodology may have a wider application, especially in medical, health care, banking, military and supply chain scenarios.

Third, data hiding techniques have been the dominated methods for protecting privacy of individual information. However, those algorithms do not pay full attention to data mining results, which may lead to sensitive rules leakages. While some algorithms are designed for preserving the rule such as with sensitive information, it may degrade the accuracy of other non-sensitive rules. Thus, further investigation, focusing on combining data and rule hiding, may be beneficial, specifically, when taking into account the interactive impact of sensitive and non-sensitive rules.

Fourth, although many machine learning methods have been used for classification, clustering, and other data mining tasks (e.g., diagnose, prediction, optimization), currently only the association rules method has been predominately used for classification. It would be interesting to see how to extend the current technique and practice into other problem domains or data mining tasks. Furthermore, it is important to find the privacy preserving technique that is independent of data mining task so that after applying privacy preserving technique a database can be released without being constrained to the original task.

Finally, identifying suitable evaluation criteria and developing benchmarks for algorithm selection are two important aspects in PPDM research. A framework for evaluating selected association rule hiding algorithms has been proposed by Bertino et al. (2005). Future research can consider testing the proposed evaluation framework for other privacy preservation algorithms, such as data distortion or cryptography methods.

#### **4. Conclusions**

PPDM has recently emerged as a new field of study. As a new comer, PPDM may offer a wide application prospect but at the same time it also brings us many issues / problems to be answered. In this study, we conduct a comprehensive survey on 29 prior studies to find out the current status of PPDM development. We propose a generic PPDM framework and a simplified taxonomy to help understand the problem and explore possible research issues. We also examine the strengths and weaknesses of different privacy preserving techniques and summarize general principles from early research to guide the selection of PPDM algorithms. As part of future work, we plan to apply the proposed evaluation framework to formally test a complete spectrum of PPDM algorithms.

#### **References**

- Bertino, E., Fovino, I., and Provenza, L. "A Framework for Evaluating Privacy Preserving Data Mining Algorithms," *Data Mining and Knowledge Discovery* (11: 2), September 2005, pp. 121-154.
- Du, W., and Zhan, Z. "Building Decision Tree Classifier on Private Data," In *Proceedings of the IEEE ICDM Workshop on Privacy, Security and Data Mining (PSDM'02)*, Maebashi City, Japan, December 2002, pp.1-8.
- Du, W., and Zhan, Z. "Using Randomized Response Techniques for Privacy-Preserving Data Mining," In *Proceedings of the Ninth ACM SIGKDD International Conference On*

- Knowledge Discovery and Data Mining*, Washington, D.C., August 24-27, 2003.
- Evfimievski, A., Gehrke, J., and Srikant, R. "Limiting Privacy Breaches in Privacy Preserving Data Mining," In *Proceedings of the Twenty-second ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems*, San Diego, California, June 09-11, 2003, pp.211-222.
- Evfimievski, A., Srikant, R., Agarwal, R., and Gehrke, J. "Privacy Preserving Mining of Association Rules," In *Proceedings of the 8th ACM SIGKDD International Conference on Knowledge Discovery in Databases and Data Mining (KDD'02)*, Edmonton, Alberta, Canada, July 2004, pp. 217-228.
- Islam, M. Z., and Brankovic, L "A Framework for Privacy Preserving Classification in Data Mining," In *Proceedings of the 2nd workshop on Australasian information security, Data Mining and Web Intelligence, and Software Internationalization (AISW'04, AWDM&WI'04, AWSI'04)*, Dunedin, New Zealand, January 2004, pp.163-168.
- Kantarcioglu, M., and Clifton, C. "Privacy-preserving Distributed Mining of Association Rules on Horizontally Partitioned Data," In *Proc. of ACM SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery (DMKD)*, June 2002.
- Kantarcioglu, M., and Clifton, C. "Assuring Privacy when Big Brother is Watching," In *Proceeding of the 8th ACM SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery, Privacy & Security*, 2003, pp.88-93.
- Klusch, M., Lodi, S., and Moro, G. "Distributed Clustering Based on Sampling Local Density Estimates," In *Proceedings of the 18th International Joint Conference on Artificial Intelligence (IJCAI'03)*, Acapulco, Mexico, August 2003, pp. 485-490.
- Lindell, Y., and Pinkas, B. "Privacy Preserving Data Mining," In *Advances in Cryptology – CRYPTO 2000*, Springer-Verlag, Aug. 20–24 2000, pp. 36–54.
- Merugu, S., and Ghosh, J. "Privacy-Preserving Distributed Clustering Using Generative Models," In *Proceedings of the 3rd IEEE International Conference on Data Mining (ICDM'03)*, Melbourne, FL, USA, November 2003, pp.211-219.
- Natwichai1, J., Li, X., and Orlowska, M. "Hiding Classification Rules for Data Sharing With Privacy Preservation," In *Proceedings of 7th International Conference on Data Warehousing and Knowledge Discovery (DaWaK'05)*, Copenhagen, Denmark, August 2005, pp. 468-477.
- Oliveira, S. R. M., and Zaïane, O. R. "Privacy Preserving Frequent Itemset Mining," In *Proceedings of the IEEE international conference on Privacy, Security and Data Mining (PSDM'02)*, Maebashi City, Japan, December 2002, pp. 43-54.
- Oliveira, S. R. M., and Zaïane, O. R. "Algorithms for Balancing Privacy and Knowledge Discovery in Association Rule Mining," In *Proceedings of the 7th International Database Engineering and Applications Symposium (IDEAS'03)*, Hong Kong, China, July 2003, pp. 54-65.
- Oliveira, S. R. M., and Zaïane, O. R. "Protecting Sensitive Knowledge By Data Sanitization," In *Proc. of the 3rd IEEE International Conference on Data Mining (ICDM'03)*, Melbourne, Florida, USA, November 2003b, pp. 613–616.
- Oliveira and, S. R. M., and Zaïane, O. R. "Privacy Preserving Clustering by Data Transformation," In *Proceedings of the 18th Brazilian Symposium on Databases*, Manaus, Amazonas, Brazil, October 2003c, pp. 304-318.
- Oliveira, S. R. M., and Zaïane, O. R. "Achieving Privacy Preservation When Sharing Data for Clustering," In *Proceedings of the International Workshop on Secure Data Management in a Connected World (SDM'04)*, In Conjunction with the 30th Very Large Data Base Conference (VLDB'04), Toronto, Canada, August 2004a, pp. 76-82.



- Oliveira, S. R. M., Zañane, O. R., and Saygin, Y. "Secure Association Rule Sharing," *PAKDD 2004b*, pp. 74-85
- Rizvi, J., and Haritsa, R. "Maintaining Data Privacy in Association Rule Mining," In *Proceedings of the 28th Very Large Data Base Conference (VLDB'02)*, Hong Kong, China, August 2002, pp. 682-693.
- Saygin, Y., Verykios, V., and Clifton, C. "Using Un-knowns to Prevent Discovery of Association Rules," *ACM SIGMOD Record* (30: 4), 2001.
- Saygin, Y., Verykios, V., and Elmagarmid, A. "Privacy Preserving Association Rule Mining", in *Proceedings of 12th Intl. Workshop on Research Issues in Data Engineering (RIDE)*, February 2002.
- Vaidya, J., and Clifton, C. "Privacy Preserving Association Rule Mining in Vertically Partitioned Data," In *Proceedings of the 8th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2002, pp. 639-644.
- Vaidya, J., and Clifton, C. "Privacy-Preserving K-Means Clustering over Vertically Partitioned Data," In *Proceedings of the 9th ACM SIGKDD International Conference on Knowledge Discovery in Data (KDD'03)*, Washington D.C., USA, August 2003, pp. 206-215.
- Vaidya, J., and Clifton, C. "Privacy-Preserving Decision Trees over Vertically Partitioned Data," In *Proceeding of the 19th Annual IFIP WG 11.3 Working Conference on Data and Applications Security (DAS'05)*, Storrs, CT, USA, August 2005, pp. 139-152.
- Veloso, A., Meira, Jr., W., Parthasarathy, S., and de Carvalho, M. "Efficient, Accurate and Privacy-Preserving Data Mining for Frequent Itemsets in Distributed Databases," In *Proceedings of the 18th Brazilian Symposium on Databases*, Manaus, Amazonas, Brazil. October 2003. pp. 281-292.
- Verykios, S., Bertino, E., Fovino, I., Provenza, L., Saygin, Y., and Theodoridis, Y. "State-of-the-art in Privacy Preserving Data Mining," *ACM SIGMOD Record* (33: 1), March 2004, pp. 50-57.
- Wang, K., Yu, S., and Chakraborty, S. "Bottom-Up Generalization: A Data Mining Solution to Privacy Protection," In *Proceedings the 4th IEEE International Conference on Data Mining (ICDM'04)*, Brighton, United Kingdom, November 2004, pp.249-256.
- Xia, Y., Yang, Y., Chi, Y., and Muntz, R. R. "Mining Association Rules with Non-uniform Privacy Concerns," Technical Report CSD-TR No. 040015, Univ. of California, 2004.
- Yang, Z., Zhong, S., Wright, R. N. "GrC. Privacy-Preserving Model Selection," In *Proceedings of the IEEE International Conference on Granular Computing*, 2006.
- Zhan, J. Z., Matwin, S., and Chang, L. "Privacy-preserving Collaborative Association Rule Mining," *DBSec 2005*, pp. 153-165.